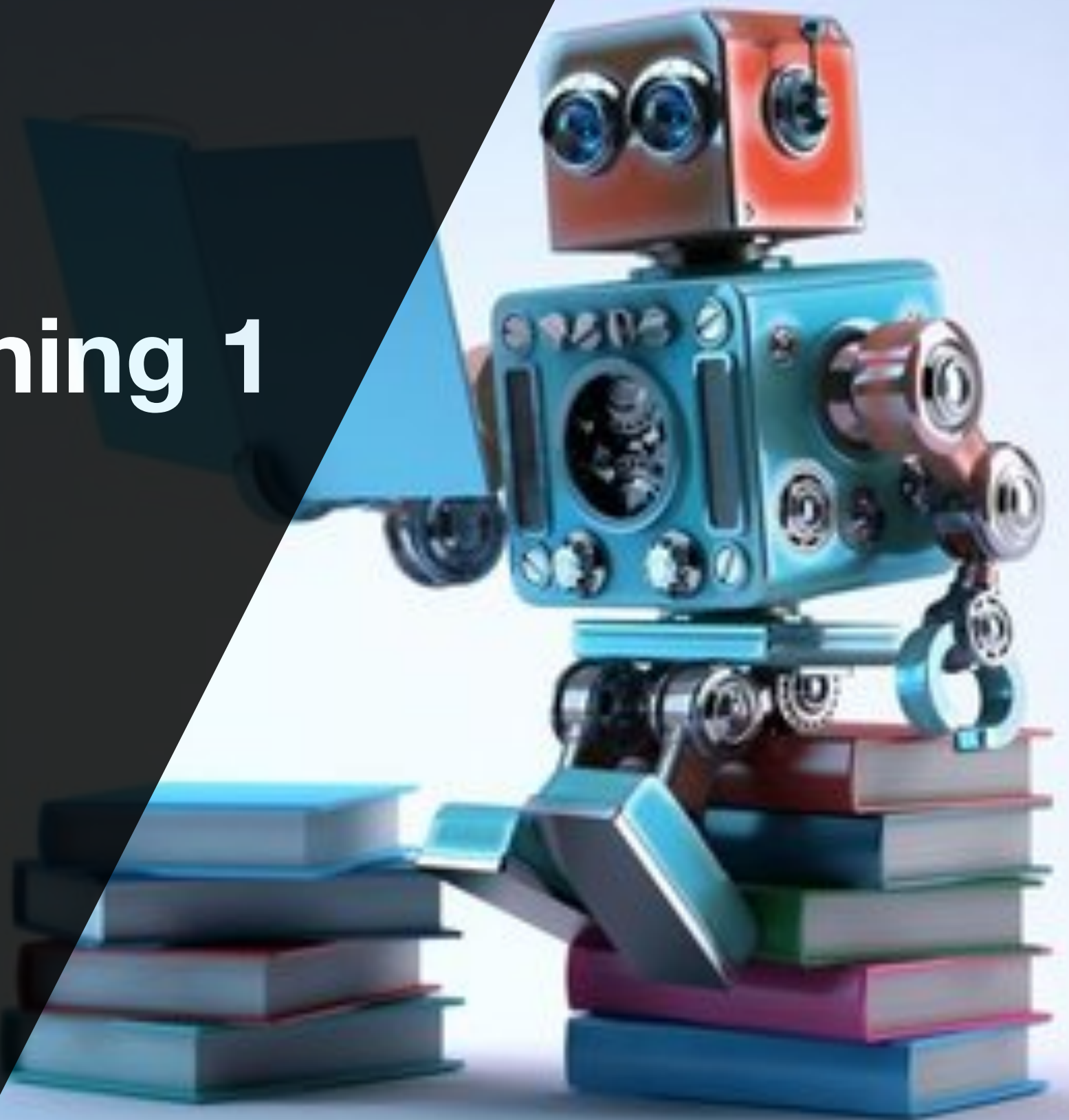


Machine Learning 1

Lecture 11.1 - Kernel Methods
Kernelizing Linear Models

Erik Bekkers

(Bishop 6.0, 6.1)

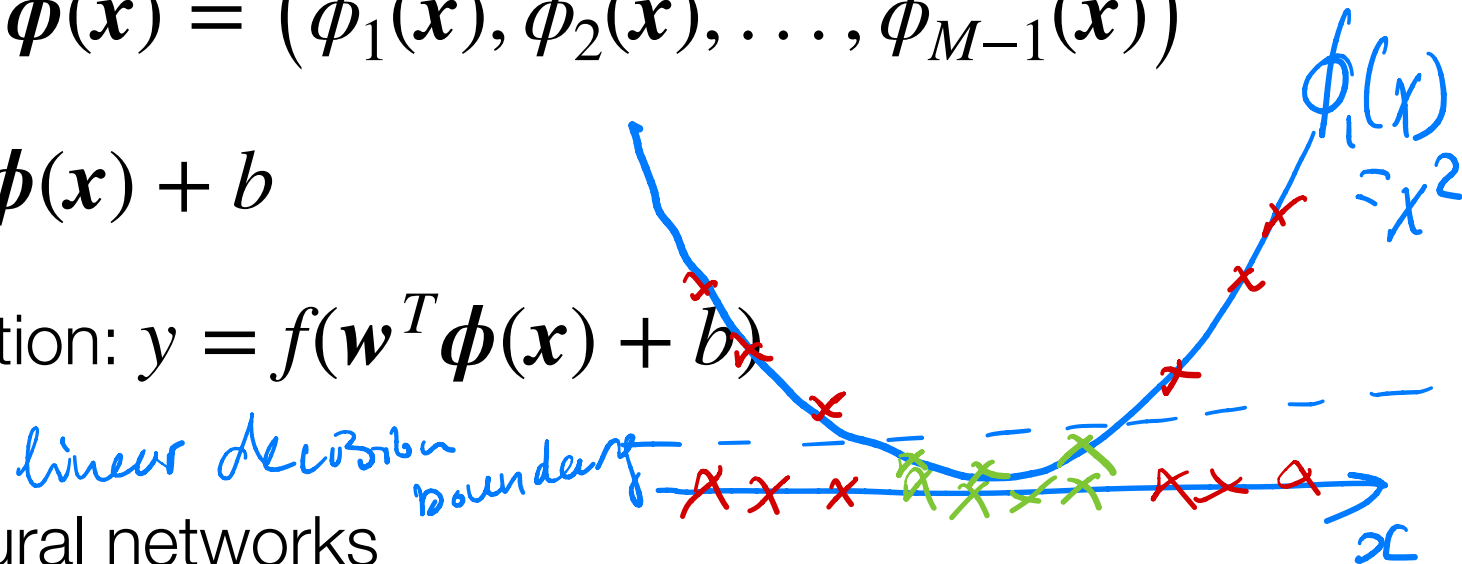


*Slide credits: Rianne van den Berg
and Patrick Forré*

Image credit: Kirillm | Getty Images

So Far: Parametric Models

- Fixed basis function methods: $\boldsymbol{\phi}(\mathbf{x}) = (\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x}))^T$
- Linear regression: $y = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b$
- Linear models for classification: $y = f(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b)$
- Learnable basis functions: Neural networks



$$y(\mathbf{x}, \mathbf{W}^{(1)}, \mathbf{w}^{(2)}) = \sum_{m=0}^M w_m^{(2)} h\left(\underbrace{\sum_{d=0}^D w_{md}^{(1)} x_d}_{\phi_m(\mathbf{x})}\right)$$

- Training:
 - MLE, MAP: use training data to obtain point estimate of \mathbf{w}
 - Full Bayesian: use training data to obtain posterior $p(\mathbf{w} | \mathbf{X}, \mathbf{t})$
- Test time: Discard training data, only need \mathbf{w} or $p(\mathbf{w} | \mathbf{X}, \mathbf{t})$

Parametric vs Non-Parametric Models

- ▶ Parametric models = models with a finite number of parameters
- ▶ Non-Parametric models = models with no explicitly defined parameters (but *implicitly still* work with (finite) or infinite number of parameters)
- ▶ Parametric methods:
 - ▶ Working in the (finite dimensional) parameter space
- ▶ Non-Parametric methods
 - ▶ Directly working in possibly infinite dimensional function spaces
 - ▶ Typically we have $M \gg N$

Non-Parametric Kernel Methods

- ▶ Kernel methods: Use (subset) of training points for predictions (test time!). Useful if $M \gg N$

- ▶ Linear parametric models:

- ▶ Can be re-cast into equivalent 'dual representation'
- ▶ Predictions are based on linear combinations of the kernel function evaluated at training data points

- ▶ For linear models with fixed feature vectors $\phi(\mathbf{x})$ we will encounter

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

- ▶ Kernel measures similarity between \mathbf{x} and \mathbf{x}' in feature space defined by mapping $\phi(\mathbf{x})$

$$k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$$

parametrized by ω and q

Kernelized Ridge Regression

- Goal: Minimize sum of squared errors with quadratic weight penalty

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{ \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) - t_n \}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

Handwritten definition of Φ :

$$\Phi = \begin{pmatrix} \phi_0(\underline{x}_1) & \phi_1(\underline{x}_1) & \dots & \phi_{M-1}(\underline{x}_1) \\ \phi_0(\underline{x}_2) & \phi_1(\underline{x}_2) & \dots & \phi_{M-1}(\underline{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\underline{x}_N) & \phi_1(\underline{x}_N) & \dots & \phi_{M-1}(\underline{x}_N) \end{pmatrix} \in \mathbb{R}^{N \times M}$$

- Solution: Solve $\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 0$:

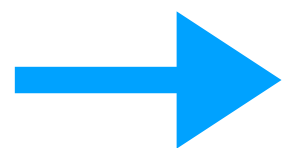
$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \sum_{n=1}^N \{ \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) - t_n \} \boldsymbol{\phi}(\mathbf{x}_n)^T + \lambda \mathbf{w}^T \mathbf{I} = 0$$

Transpose on both sides, using $(\mathbf{a}^T \mathbf{B})^T = \mathbf{B}^T \mathbf{a}$

$$\Leftrightarrow \mathbf{w}^T \left(\sum_{n=1}^N \boldsymbol{\phi}(\mathbf{x}_n) \boldsymbol{\phi}(\mathbf{x}_n)^T + \lambda \mathbf{I} \right) = \sum_{n=1}^N t_n \boldsymbol{\phi}(\mathbf{x}_n)^T$$

$\mathbf{A} \mathbf{x} = \mathbf{b} \Rightarrow \mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$

$$\Leftrightarrow \left(\sum_{n=1}^N \boldsymbol{\phi}(\mathbf{x}_n) \boldsymbol{\phi}(\mathbf{x}_n)^T + \lambda \mathbf{I} \right) \mathbf{w} = \sum_{n=1}^N t_n \boldsymbol{\phi}(\mathbf{x}_n)$$



$$\mathbf{w} = \left(\sum_{n=1}^N \boldsymbol{\phi}(\mathbf{x}_n) \boldsymbol{\phi}(\mathbf{x}_n)^T + \lambda \mathbf{I} \right)^{-1} \sum_{n=1}^N t_n \boldsymbol{\phi}(\mathbf{x}_n) = \underbrace{(\Phi^T \Phi + \lambda \mathbf{I})^{-1}}_{M \times M} \Phi^T \mathbf{t}$$

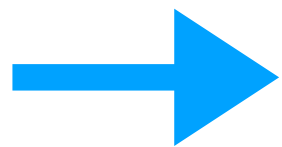
$M \times M$

Kernelized Ridge Regression

- Goal: Minimize sum of squared errors with quadratic weight penalty

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- Solution: Solve $\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 0$:



$$\mathbf{w} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I}_M)^{-1} \boldsymbol{\Phi}^T \mathbf{t},$$

$$\boldsymbol{\Phi}^T \boldsymbol{\Phi} \in \mathbb{R}^{M \times M}$$

- Use matrix inversion lemma (see e.g. Bishop C.5):

$$\left(P^{-1} + B^T R^{-1} B \right)^{-1} B^T R^{-1} = P B^T \left(B P B^T + R \right)^{-1}$$

$\left(\lambda \mathbf{I}_M + \underbrace{\boldsymbol{\Phi}^T \mathbf{I}_N \boldsymbol{\Phi}}_{\boldsymbol{\Phi} \boldsymbol{\Phi}^T} \right)^{-1} \boldsymbol{\Phi}^T$

$$\begin{pmatrix} P^{-1} = \lambda \mathbf{I}_M \\ B = \boldsymbol{\Phi} \\ R = \mathbf{I}_N \end{pmatrix}$$

- Allows us to alternatively obtain \mathbf{w} via

$$\mathbf{w} = \boldsymbol{\Phi}^T (\boldsymbol{\Phi} \boldsymbol{\Phi}^T + \lambda \mathbf{I}_N)^{-1} \mathbf{t} = \boldsymbol{\Phi}^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

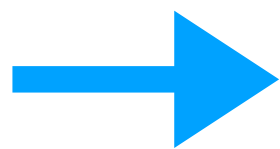
- With Gram matrix $\mathbf{K} = \boldsymbol{\Phi} \boldsymbol{\Phi}^T$ with $K_{ij} = \boldsymbol{\phi}^T(\mathbf{x}_i) \boldsymbol{\phi}(\mathbf{x}_j)$

Kernelized Ridge Regression

- ▶ Goal: Minimize sum of squared errors with quadratic weight penalty

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{ \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) - t_n \}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- ▶ Solution: Solve $\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 0$:



$$\mathbf{w} = \underbrace{\Phi^T (K + \lambda \mathbf{I}_N)^{-1} \mathbf{t}}_{\mathbf{a}}$$

$$\mathbf{K} = \Phi \Phi^T, \quad K_{ij} = \boldsymbol{\phi}^T(\mathbf{x}_i) \boldsymbol{\phi}(\mathbf{x}_j)$$

- ▶ Primal/dual viewpoint

- ▶ Primal variable: $\mathbf{w} = \Phi^T \mathbf{a}$

$$\arg \min_{\mathbf{w}, \mathbf{z}} \frac{1}{2} \|\mathbf{z}\|^2 \text{ with constraints } \mathbf{z} = \Phi \mathbf{w} - \mathbf{t} \text{ and } \frac{1}{2} \|\mathbf{w}\|^2 \leq R^2$$

- ▶ Dual variable: $\mathbf{a} = (K + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$

$$\arg \min_{\mathbf{a}} \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{t} + \frac{\lambda}{2} \|\mathbf{a}\|^2$$

- ▶ Predictive mean of primal viewpoint $y(\mathbf{x}', \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}')$

$$\text{of dual viewpoint } y(\mathbf{x}', \mathbf{a}) = \sum_{n=1}^N a_n k(\mathbf{x}_n, \mathbf{x}')$$

Primal vs Dual/Kernel Approach

- Computational cost (closed form solutions):

- The dual variables. $\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$ $O(N^3)$

- The primal variables $\mathbf{w} = (\Phi^T \Phi + \lambda \mathbf{I}_M)^{-1} \Phi^T \mathbf{t}$ $O(M^3)$

- Computational cost (predictions):

- Dual case: $y(\mathbf{x}', \mathbf{a}) = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x}')$ $O(NM)$

- Primal case: $y(\mathbf{x}', \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x}')$ $O(M)$

- But... dual approach:

- No explicit parameters (implicitly many!) -> nonparametric model

- Does not rely on explicit features but on similarity kernel function.

- Can be slow at prediction

- Upcoming: **Kernel methods with sparse solutions!**

$N' \ll N$
 $O(N' M)$